

from:
N.J., Smelter, & P.B., Balthes (Eds.) (2001).
Encyclopedia of the Social and Behavioral Sciences.
London: Elsevier Science.

Article Title: *Linear Algebra for Neural Networks*

By: *Hervé Abdi*

Author Address: Hervé Abdi, School of Human Development, MS: Gr.4.1,
The University of Texas at Dallas, Richardson, TX 750833-0688, USA

Phone: 972 883 2065, **fax:** 972 883 2491 **Date:** June 1, 2001

E-mail: herve@utdallas.edu

Abstract

Neural networks are quantitative models which learn to associate input and output patterns adaptively with the use of learning algorithms. We expose four main concepts from linear algebra which are essential for analyzing these models: 1) the projection of a vector, 2) the eigen and singular value decomposition, 3) the gradient vector and Hessian matrix of a vector function, and 4) the Taylor expansion of a vector function. We illustrate these concepts by the analysis of the Hebbian and Widrow-Hoff rules and some basic neural network architectures (i.e., the linear autoassociator, the linear heteroassociator, and the error backpropagation network). We show also that neural networks are equivalent to iterative versions of standard statistical and optimization models such as multiple regression analysis and principal component analysis.

1 Introduction

Linear algebra is particularly well suited to analyze the class of neural networks called *associators*. These quantitative models learn to associate input and out-

put patterns adaptively with the use of learning algorithms. When the set of input patterns is different from the set of output patterns, the models are called *heteroassociators*. When input patterns and output patterns are the same, the models are called *autoassociators*. Associators consist of layers of elementary units called neurons. Information flows from the first layer (the input layer), to the last one (the output layer). Some architectures may include intermediate layers (the hidden layers). Typically neurons from a given layer are connected to the neurons of another layer. Linear algebra operations describe the transformations of the information as it flows from one layer to another one.

2 Prerequisite notions and notations

Vectors are represented by lower case boldface letters (*e.g.*, \mathbf{x}), matrices by upper case boldface letters (*e.g.*, \mathbf{X}). The following elementary notions are supposed to be known (see Abdi, Valentin & Edelman, 1999, for a neural network approach): the transpose operation, (*e.g.*, \mathbf{x}^\top), the norm of a vector (*e.g.*, $\|\mathbf{x}\|$), the scalar product (*e.g.*, $\mathbf{x}^\top \mathbf{w}$) and the product of two matrices (*e.g.*, $\mathbf{A}\mathbf{B}$). We will use also the elementwise or Hadamar product (*e.g.*, $\mathbf{A} \otimes \mathbf{B}$).

3 Projection

3.1 Cosine between two vectors

The *cosine* of vectors \mathbf{x} and \mathbf{y} is the cosine of the angle made by the origin of the space and the points defined by the coordinates of the vectors. Thus,

$$\cos(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^\top \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} . \quad (1)$$

The cosine reflects the *similarity* between vectors. When two vectors are proportional to each other (*i.e.*, when they have the same direction), their cosine is equal to 1; when they are orthogonal to each other, their cosine is equal to 0.

3.2 Distance between vectors

Among the large family of *distances* between vectors, the most popular by far is the *Euclidean* distance. It is closely related to the cosine between vectors and is defined as

$$\begin{aligned} d^2(\mathbf{x}, \mathbf{y}) &= (\mathbf{x} - \mathbf{y})^\top (\mathbf{x} - \mathbf{y}) = \mathbf{x}^\top \mathbf{x} + \mathbf{y}^\top \mathbf{y} - 2\mathbf{x}^\top \mathbf{y} \\ &= \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - 2[\cos(\mathbf{x}, \mathbf{y}) \times \|\mathbf{x}\| \times \|\mathbf{y}\|] = \sum (x_i - y_i)^2 . \end{aligned} \quad (2)$$

3.3 Projection of one vector onto another vector

The (orthogonal) projection of vector \mathbf{x} on vector \mathbf{w} is defined as

$$\text{proj}_{\mathbf{w}}\mathbf{x} = \frac{\mathbf{x}^T\mathbf{w}}{\mathbf{w}^T\mathbf{w}}\mathbf{w} = \cos(\mathbf{x}, \mathbf{w}) \times \frac{\|\mathbf{x}\|}{\|\mathbf{w}\|}\mathbf{w} . \quad (3)$$

The norm of $\text{proj}_{\mathbf{w}}\mathbf{x}$ is its distance to the origin of the space. It is equal to

$$\|\text{proj}_{\mathbf{w}}\mathbf{x}\| = \frac{|\mathbf{x}^T\mathbf{w}|}{\|\mathbf{w}\|} = |\cos(\mathbf{x}, \mathbf{w})| \times \|\mathbf{x}\| . \quad (4)$$

3.4 Illustration: Hebbian and Widrow-Hoff learning rules

A neural network consists of cells connected to other cells via modifiable weighted *connections* called *synapses*. Consider a network made of I input cells and only one output cell. The information is transmitted via the synapses, from a set of external input cells to the output cell which gives a response corresponding to its state of activation. If the input pattern and the set of synaptic weights are given by I -dimensional vectors denoted \mathbf{x} , and \mathbf{w} , the activation of the output cell is obtained as

$$a = \mathbf{x}^T\mathbf{w} . \quad (5)$$

So the activation is proportional to the norm of the projection of the input vector onto the weight vector. The *response* or *output* of the cell is denoted o . For a *linear cell*, it is proportional to the activation (for convenience, assume that the proportionality constant is equal to 1). *Linear heteroassociators* and *autoassociators* are made of linear cells. In general, the output of a cell is a *function* (often, but not necessarily, continuous), called the *transfer function*, of its activation

$$o = f(a) . \quad (6)$$

For example, in *backpropagation networks*, the (nonlinear) transfer function is usually the logistic function

$$o = f(a) = \text{logist } \mathbf{w}^T\mathbf{x} = \frac{1}{1 + \exp\{-a\}} . \quad (7)$$

Often, a neural network is designed to associate, to a given input, a specific response called the target, denoted t . *Learning* is equivalent to defining a rule that specifies how to add a small quantity to each synaptic weight at each learning iteration. Learning makes the output of the network closer to the target.

Learning rules come in two main flavors: *supervised* (e.g., Widrow-Hoff) which take into account the error or distance between the response of the neuron and the target and *unsupervised* (e.g., Hebbian) which do not require such “feedback.” The Hebbian learning rule modifies the weight vector at iteration $n + 1$ as

$$\mathbf{w}_{[n+1]} = \mathbf{w}_{[n]} + \eta t\mathbf{x} , \quad (8)$$

where η is a small positive constant called the *learning constant*. So, a Hebbian learning iteration moves the weight vector in the direction of the input vector by an amount proportional to the target.

The Widrow-Hoff learning rule uses the error and the derivative of the transfer function (f') to compute the correction as

$$\mathbf{w}_{[n+1]} = \mathbf{w}_{[n]} + \eta f'(a)(t - o)\mathbf{x} . \quad (9)$$

So, a Widrow-Hoff learning iteration moves the weight vector in the direction of the input vector by an amount proportional to the error.

For networks with several cells (say J) in the output layer, the pattern of activation, output, and target become J -dimensional vectors (denoted \mathbf{a} , \mathbf{o} , and \mathbf{t} , respectively), and the synaptic weights are stored in an $I \times J$ matrix \mathbf{W} . The learning equations become

$$\begin{aligned} \mathbf{W}_{[n+1]} &= \mathbf{W}_{[n]} + \eta \mathbf{x} \mathbf{t}^\top \text{ (Hebbian) and} \\ \mathbf{W}_{[n+1]} &= \mathbf{W}_{[n]} + \eta (f'(\mathbf{a}) \otimes \mathbf{x})(\mathbf{t} - \mathbf{o})^\top \text{ (Widrow-Hoff) ,} \end{aligned} \quad (10)$$

(the derivative of the transfer function is applied elementwise).

In general, several (say K) input/target associations are to be learned. Then the set of input patterns is stored in an $I \times K$ matrix denoted \mathbf{X} , the activation and target patterns respectively are stored in $J \times K$ matrices denoted \mathbf{A} and \mathbf{T} respectively. The activation and learning iterations can be computed for all the patterns at once (this is called *batch* mode). The output matrix is computed as

$$\mathbf{O} = f(\mathbf{A}) = f(\mathbf{W}\mathbf{X}^\top) \quad (11)$$

(f is applied elementwise). The learning equations become

$$\mathbf{W}_{[n+1]} = \mathbf{W}_{[n]} + \eta \mathbf{X} \mathbf{T}^\top \text{ (Hebb) and} \quad (12)$$

$$\mathbf{W}_{[n+1]} = \mathbf{W}_{[n]} + \eta (f'(\mathbf{A}) \otimes \mathbf{X})(\mathbf{T} - \mathbf{O})^\top \text{ (Widrow-Hoff) .} \quad (13)$$

4 Eigenvalues, eigenvectors, and the singular value decomposition

Eigenvectors of a given square matrix \mathbf{W} (resulting from its *eigendecomposition*) are vectors invariant under multiplication by \mathbf{W} . The eigendecomposition is best defined for a subclass of matrices called *positive semi-definite* matrices. A matrix \mathbf{X} is positive semi-definite if there exists another matrix \mathbf{Y} such that $\mathbf{X} = \mathbf{Y}\mathbf{Y}^\top$. This is the case for most matrices used in neural networks and so we will consider only this case here.

Formally, a (nonzero) vector \mathbf{u} is an eigenvector of a square matrix \mathbf{W} if

$$\lambda \mathbf{u} = \mathbf{W} \mathbf{u} . \quad (14)$$

The scalar λ is the *eigenvalue* associated with \mathbf{u} . So, \mathbf{u} is an eigenvector of \mathbf{W} if its direction is invariant under multiplication by \mathbf{W} (only its length is changed

if $\lambda \neq 1$). There are, in general, several eigenvectors for a given matrix (at most as many as the dimension of \mathbf{W}). They are in general ordered by decreasing order of their eigenvalue. So, the first eigenvector, \mathbf{u}_1 has the largest eigenvalue λ_1 . The number of eigenvectors with a non-zero eigenvalue is the *rank* of the matrix.

The eigenvalues of positive semi-definite matrices are always positive or zero (a matrix with strictly positive eigenvalues, is *positive definite*). Also, any two eigenvectors with different eigenvalues are orthogonal, i.e.

$$\mathbf{u}_\ell^\top \mathbf{u}_{\ell'} = 0 \quad \forall \ell \neq \ell' . \quad (15)$$

In addition, the set of eigenvectors of a matrix constitutes an orthogonal basis for its rows and columns. This is expressed by defining two matrices: the eigenvector matrix \mathbf{U} , and the diagonal matrix of the eigenvalues: $\mathbf{\Lambda}$. The eigendecomposition of \mathbf{W} (with rank L) is

$$\mathbf{W} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top = \sum_{\ell}^L \lambda_\ell \mathbf{u}_\ell \mathbf{u}_\ell^\top, \text{ or equivalently: } \mathbf{\Lambda} = \mathbf{U}^\top \mathbf{W} \mathbf{U} . \quad (16)$$

The *singular value decomposition* (SVD) generalizes the eigendecomposition to rectangular matrices. If \mathbf{X} is an $I \times K$ matrix, its SVD is defined as

$$\mathbf{X} = \mathbf{U}\mathbf{\Delta}\mathbf{V}^\top \text{ with } \mathbf{U}^\top \mathbf{U} = \mathbf{V}^\top \mathbf{V} = \mathbf{I} \text{ and } \mathbf{\Delta} \text{ being a diagonal matrix} . \quad (17)$$

(\mathbf{I} being the *identity* matrix). The diagonal elements of $\mathbf{\Delta}$ are real positive numbers called the *singular values* of \mathbf{X} . The matrices \mathbf{U} and \mathbf{V} are the left and right matrices of singular vectors (which are also eigenvectors, see below). The SVD is closely related to the eigendecomposition because \mathbf{U} , \mathbf{V} , and $\mathbf{\Delta}$ can be obtained from the eigendecomposition of matrices $\mathbf{X}^\top \mathbf{X}$ and $\mathbf{X} \mathbf{X}^\top$ as

$$\mathbf{X} \mathbf{X}^\top = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top, \quad \mathbf{X}^\top \mathbf{X} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top, \text{ and } \mathbf{\Delta} = \mathbf{\Lambda}^{\frac{1}{2}} \quad (18)$$

(note that $\mathbf{X}^\top \mathbf{X}$ and $\mathbf{X} \mathbf{X}^\top$ have the same eigenvalues).

The eigen- and singular value decompositions are used in most fields of applied mathematics including statistics, image processing, mechanics, and dynamical systems. For neural networks, they are essential for studying the dynamics of linear autoassociators and heteroassociators.

4.1 Iterative processes

A linear heteroassociator using the Widrow-Hoff rule, learning modifies only the eigenvalues of the weight matrix. Specifically, if the patterns to be learned are stored in an $I \times K$ matrix \mathbf{X} , with singular value decomposition as $\mathbf{X} = \mathbf{U}\mathbf{\Delta}\mathbf{V}^\top$, then the learning equation (see Equation 12) becomes [because for a linear heteroassociator, $\mathbf{O} = \mathbf{A}$, and $f'(\mathbf{A}) = \mathbf{I}$]

$$\mathbf{W}_{[n+1]} = \mathbf{W}_{[n]} + \eta \mathbf{X} (\mathbf{T} - \mathbf{A})^\top = \mathbf{U} \left\{ \mathbf{\Delta}^{-1} \left[\mathbf{I} - (\mathbf{I} - \eta \mathbf{\Delta}^2)^{n+1} \right] \right\} \mathbf{V}^\top \mathbf{T}^\top, \quad (19)$$

(see Abdi, 1994, p.54ff.).

The Hebbian weight matrix corresponds to the first iteration of the algorithm, i.e.,

$$\mathbf{W}_{[1]} = \mathbf{U} \left\{ \Delta^{-1} \left[\mathbf{I} - (\mathbf{I} - \eta \Delta^2)^1 \right] \right\} \mathbf{V}^T \mathbf{T}^T = \mathbf{U} \{ \eta \Delta \} \mathbf{V}^T \mathbf{T}^T = \eta \mathbf{X} \mathbf{T}^T . \quad (20)$$

Equation 19 characterizes the values of η that allow the iterative process to converge. Denoting by δ_{\max} the largest singular value of \mathbf{X} , if η is such that

$$0 < \eta < 2\delta_{\max}^{-2} \quad (21)$$

then it can be shown that (see Abdi, 1994)

$$\lim_{n \rightarrow \infty} (\mathbf{I} - \eta \Delta^2)^n = \mathbf{0}, \text{ and } \lim_{n \rightarrow \infty} \mathbf{W}_{[n]} = \mathbf{U} \Delta^{-1} \mathbf{V}^T \mathbf{T}^T = \mathbf{X}^+ \mathbf{T} . \quad (22)$$

The matrix $\mathbf{X}^+ = \mathbf{U} \Delta^{-1} \mathbf{V}^T$ is the *pseudo-inverse* of \mathbf{X} . It gives a least-square optimal solution for the association between the input and the target. Thus, the linear heteroassociator is equivalent to *linear multiple regression*. If η is outside the interval defined by Equation 21, then both the singular values and the elements of the weight matrix will grow at each iteration. In practice, because neural networks are simulated by digital computers, the weight matrix will eventually reach the limits of the precision of the machine.

When the target vectors are the same as the input vectors (*i.e.*, when each input is associated to itself), the linear heteroassociator becomes a linear autoassociator. The previous approach shows that, now, the Hebbian weight matrix is the cross-product matrix

$$\mathbf{W}_{[1]} = \mathbf{X} \mathbf{X}^T = \mathbf{U} \Lambda \mathbf{U}^T . \quad (23)$$

With Widrow-Hoff learning, when convergence is reached, all nonzero eigenvalues of the weight matrix are equal to 1. The weight matrix is then said to be *spherical*; it is equal to

$$\mathbf{W}_{[\infty]} = \mathbf{U} \mathbf{U}^T . \quad (24)$$

Because the statistical technique of *principal component analysis* (PCA) computes the eigendecomposition of a cross-product matrix similar to \mathbf{W} , the linear autoassociator is considered as the neural network equivalent of PCA.

5 Optimization, Derivative and Matrices

Neural networks are often used to optimize a function of the synaptic weights. The *differentiation* of a function is the major concept for exploring *optimization* problems and, for neural networks, it involves differentiating vectors or matrix functions. In this context, we need to consider the transfer function as being a function of the weight vector. This is expressed by rewriting Equation 6 as

$$o = f(\mathbf{w}) . \quad (25)$$

The derivative of $f(\mathbf{w})$ with respect to the $I \times 1$ vector \mathbf{w} is denoted by $\nabla_{f(\mathbf{w})}$. It is also called the *gradient* of f , i.e.,

$$\nabla_{f(\mathbf{w})} = \frac{\partial f}{\partial \mathbf{w}} = \left[\frac{\partial f}{\partial w_1}, \dots, \frac{\partial f}{\partial w_i}, \dots, \frac{\partial f}{\partial w_I} \right]^\top. \quad (26)$$

For example, the derivative of the output of a linear neuron is

$$\frac{\partial f}{\partial \mathbf{w}} = \left[\frac{\partial \mathbf{w}^\top \mathbf{x}}{\partial w_1}, \dots, \frac{\partial \mathbf{w}^\top \mathbf{x}}{\partial w_i}, \dots, \frac{\partial \mathbf{w}^\top \mathbf{x}}{\partial w_I} \right]^\top = [x_1, \dots, x_i, \dots, x_I]^\top = \mathbf{x}. \quad (27)$$

When a function is twice differentiable, the second order derivatives are stored in a matrix called the *Hessian* matrix of the function. It is often denoted by \mathbf{H} or ∇_f^2 and is formally defined as

$$\mathbf{H} = \nabla_f^2 = \begin{bmatrix} \frac{\partial^2 f}{\partial w_1^2} & \frac{\partial^2 f}{\partial w_1 w_2} & \cdots & \frac{\partial^2 f}{\partial w_1 w_I} \\ \frac{\partial^2 f}{\partial w_2 w_1} & \frac{\partial^2 f}{\partial w_2^2} & \cdots & \frac{\partial^2 f}{\partial w_2 w_I} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial w_I w_1} & \frac{\partial^2 f}{\partial w_I w_2} & \cdots & \frac{\partial^2 f}{\partial w_I^2} \end{bmatrix}. \quad (28)$$

5.1 Conditions for minimum

A standard problem is to show that a given learning rule finds an optimum solution in the sense that a function of the weight vector (or matrix) called the *error function* reaches its minimum value when learning has converged. Often, the error function is defined as the sum of the squared error over all patterns.

When the gradient of the error function can be evaluated, a necessary condition for optimality (i.e., either minimum or maximum) is to find a weight vector $\tilde{\mathbf{w}}$ such that

$$\nabla_{f(\tilde{\mathbf{w}})} = \mathbf{0}. \quad (29)$$

This condition is also sufficient provided \mathbf{H} is positive definite (cf. Haykin, 1999).

5.2 Taylor expansion

The Taylor expansion is the standard technique used to obtain a linear or a quadratic approximation of a function of one variable. Recall that the Taylor expansion of a continuous function $f(x)$ is

$$\begin{aligned} f(x) &= f(a) + (x-a) \frac{f'(a)}{1!} + (x-a)^2 \frac{f''(a)}{2!} + \dots + (x-a)^n \frac{f^{[n]}(a)}{n!} + \dots \\ &= f(a) + (x-a) \frac{f'(a)}{1!} + (x-a)^2 \frac{f''(a)}{2!} + \mathcal{R}_2. \end{aligned} \quad (30)$$

(where \mathcal{R}_2 represents all the terms of higher order than 2, and a is a “convenient” value at which to evaluate f).

This technique can be extended to matrix and vector functions. It involves the notion of gradient and Hessian. Now a vector function $f(\mathbf{x})$ is expressed as:

$$f(\mathbf{x}) = f(\mathbf{a}) + f(\mathbf{x} - \mathbf{a})^\top \nabla_{f(\mathbf{a})} + f(\mathbf{x} - \mathbf{a})^\top \nabla_{f(\mathbf{a})}^2 f(\mathbf{x} - \mathbf{a}) + \mathcal{R}_2 . \quad (31)$$

5.3 Iterative minimization

A learning rule can be shown to converge to an optimum if it diminishes the value of the error function at each iteration. When the gradient of the error function can be evaluated, the *gradient* technique (or *steepest descent*) adjusts the weight vector by moving it in the direction opposite to the gradient of the error function. Formally, the correction for the $(n + 1)$ -th iteration is

$$\mathbf{w}_{[n+1]} = \mathbf{w}_{[n]} + \Delta = \mathbf{w}_{[n]} - \eta \nabla_{f(\mathbf{w})} \quad (32)$$

(where $\nabla_{f(\mathbf{w})}$ is computed for $\mathbf{w}_{[n]}$).

As an example, let us show that for a linear heteroassociator, the Widrow-Hoff learning rule minimizes iteratively the squared error between target and output. The error function is

$$e^2 = (t - o)^2 = t^2 + o^2 - 2to = t^2 + \mathbf{x}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x} - 2t \mathbf{w}^\top \mathbf{x} . \quad (33)$$

The gradient of the error function is

$$\frac{\partial e}{\partial \mathbf{w}} = 2(\mathbf{w}^\top \mathbf{x}) \mathbf{x} - 2t \mathbf{x} = -2(t - \mathbf{w}^\top \mathbf{x}) \mathbf{x} . \quad (34)$$

The weight vector is corrected by moving it in the opposite direction of the gradient. This is obtained by adding a small vector denoted $\Delta_{\mathbf{w}}$ opposite to the gradient. This gives the following correction for iteration $n + 1$:

$$\mathbf{w}_{[n+1]} = \mathbf{w}_{[n]} + \Delta_{\mathbf{w}} = \mathbf{w}_{[n]} - \eta \frac{\partial e}{\partial \mathbf{w}} = \mathbf{w}_{[n]} + \eta(t - \mathbf{w}^\top \mathbf{x}) \mathbf{x} = \mathbf{w}_{[n]} + \eta(t - o) \mathbf{x} . \quad (35)$$

This gives the rule defined by Equation 9.

The gradient method works because the gradient of $\mathbf{w}_{[n]}$ is a first order Taylor approximation of the gradient of the optimal weight vector $\tilde{\mathbf{w}}$. It is a favorite technique in neural networks because the popular error backpropagation is a gradient technique.

Newton's method is a second order Taylor approximation, it uses the inverse of the Hessian of \mathbf{w} (supposing it exists). It gives a better numerical approximation but necessitates more computation. Here the correction for iteration $n + 1$ is

$$\mathbf{w}_{[n+1]} = \mathbf{w}_{[n]} + \Delta = \mathbf{w}_{[n]} - (\mathbf{H}^{-1})(\nabla_{f(\mathbf{w})}) \quad (36)$$

(where $\nabla_{f(\mathbf{w})}$ is computed for $\mathbf{w}_{[n]}$).

6 Useful References

Linear algebra at the level of this presentation is available in the following recent books: Abdi *et al.* (1999), Bishop (1995) Ellacot and Bose (1996), Haggan, Demuth, and Beale (1996), Haykin (1999), Reed and Marks (1999), Ripley (1996), and Rojas (1996).

See also: Artificial neural networks: neurocomputation; Backpropagation; Hebb, Donald Olding (1904–1985); Statistical pattern recognition.

References

- [1] ABDI, H. (1994a) *Les réseaux de neurones*. Grenoble, France: PUG.
- [2] ABDI, H., VALENTIN, D., & EDELMAN, B. (1999) *Neural networks*. Thousand Oak, CA: Sage.
- [3] BISHOP, C.M. (1995) *Neural network for pattern recognition*. Oxford, UK: Oxford University Press.
- [4] ELLACOTT, S., & BOSE, D. (1996) *Neural networks: Deterministic methods of analysis*. London: ITC.
- [5] HAGAN, M. T., DEMUTH, H. B., & BEALE, M. (1996) *Neural networks design*. Boston: PWS.
- [6] HAYKIN, S. (1999) *Neural networks: A comprehensive foundation* (2nd ed). New York: Prentice Hall.
- [7] REED, R.D., MARKS R.J. (1999) *Neural smithing*. Cambridge, MA: MIT press.
- [8] RIPLEY, B.D. (1996) *Pattern recognition and neural networks*. Cambridge, MA: Cambridge University Press.
- [9] ROJAS, R. (1996) *Neural networks*. New York: Springer-Verlag.